

# Covariance Matrices and Influence Scores for Mean Field Variational Bayes

Ryan Giordano  
Department of Statistics  
University of California, Berkeley  
Berkeley, CA 94720  
rgiordano@berkeley.edu

Tamara Broderick  
Department of EECS,  
Massachusetts Institute of Technology  
Cambridge, MA 02139  
tbroderick@csail.mit.edu

February 27, 2015

## Abstract

Mean field variational Bayes (MFVB) is a popular posterior approximation method due to its fast runtime on large-scale data sets. However, it is well known that a major failing of MFVB is that it underestimates the uncertainty of model variables (sometimes severely) and provides no information about model variable covariance. We develop a fast, general methodology for exponential families that augments MFVB to deliver accurate uncertainty estimates for model variables—both for individual variables and coherently across variables. MFVB for exponential families defines a fixed-point equation in the means of the approximating posterior, and our approach yields a covariance estimate by perturbing this fixed point. Inspired by linear response theory, we call our method linear response variational Bayes (LRVB). We also show how LRVB can be used to quickly calculate a measure of the influence of individual data points on parameter point estimates. We demonstrate the accuracy and scalability of our method by learning Gaussian mixture models for both simulated and real data.

## 1 Introduction

With increasingly efficient data collection methods, scientists are interested in quickly analyzing ever larger data sets. In particular, the promise of these large data sets is not simply to fit old models but instead to learn more nuanced patterns from data than has been possible in the past. In theory, the Bayesian paradigm promises exactly these desiderata. Hierarchical modeling allows practitioners to capture complex relationships between variables of interest. Moreover, Bayesian analysis allows practitioners to quantify the uncertainty in any model estimates—and to do so coherently across all of the model variables.

*Mean field variational Bayes* (MFVB), a method for approximating a Bayesian posterior distribution, has grown in popularity due to its fast runtime on large-scale data sets [1–3]. But it is well known that a major failing of MFVB is that it gives underestimates of the uncertainty of model variables that can be arbitrarily bad, even when approximating a simple multivariate Gaussian distribution [4–6], and provides no information about how the uncertainties in different model variables interact [5–8]. We develop a fast, general methodology for exponential families that augments MFVB to deliver accurate uncertainty estimates for model variables—both for individual variables

and coherently across variables. In particular, as we elaborate in Section 2, MFVB for exponential families defines a fixed-point equation in the means of the approximating posterior, and our approach yields a covariance estimate by perturbing this fixed point. The perturbations of *linear response theory* have previously been applied for machine learning by [9] and specifically for mean-field methods by [10] and [11]. Our contribution is to use exponential families to derive particularly simple and scalable formulas for covariance estimation and to develop a method to quickly calculate *influence scores*, which measure the influence of individual data points on parameter point estimates. We call our method *linear response variational Bayes* (LRVB).

We demonstrate the accuracy and scalability of our LRVB covariance estimates with experiments that focus on finite mixtures of multivariate Gaussians, which have historically been a sticking point for MFVB covariance estimates [5, 6]. We employ simulated data as well as the MNIST handwritten digit data set [12]. We show that the LRVB variance estimates are nearly identical to those produced by a Markov Chain Monte Carlo (MCMC) sampler, even when MFVB variance is dramatically underestimated. For these mixture models, we show that LRVB gives accurate covariance estimates orders of magnitude faster than MCMC on a wide range of problems. We demonstrate both theoretically and empirically that, for this Gaussian mixture model, LRVB scales linearly in the number of data points and approximately quadratically in the dimension of the parameter space. Finally, we show how LRVB allows fast computation of the influence scores mentioned above.

## 2 Mean-field variational Bayes in exponential families

Denote our  $N$  observed data points by the  $N$ -long column vector  $x$ , and denote our unobserved model parameters by  $\theta$ . Here,  $\theta$  is a column vector residing in some space  $\Theta$ ; it has  $J$  subgroups and total dimension  $D$ . Our model is specified by a distribution of the observed data given the model parameters—the likelihood  $p(x|\theta)$ —and a prior distributional belief on the model parameters  $p(\theta)$ . Bayes’ Theorem yields the posterior  $p(\theta|x)$ .

MFVB approximates  $p(\theta|x)$  by a factorized distribution of the form  $q(\theta) = \prod_{j=1}^J q(\theta_j)$  such that the Kullback-Liebler divergence  $\text{KL}(q||p)$  between  $q$  and  $p$  is minimized:

$$\begin{aligned} q^* &:= \arg \min_q \text{KL}(q||p) \\ &= \arg \min_q \mathbb{E}_q \left[ \log p(\theta|x) - \sum_{j=1:J} \log q(\theta_j) \right]. \end{aligned}$$

By the assumed  $q$  factorization, the solution to this minimization obeys the following fixed point equations [5]:

$$\log q_j^*(\theta_j) = \mathbb{E}_{q_i^*: i \in [J] \setminus j} \log p(\theta, x) + C. \quad (1)$$

Here, and for the rest of the text,  $C$  denotes a constant and  $[J] := \{1, \dots, J\}$ . For index  $j$ , suppose that  $p(\theta_j|\theta_{i \in [J] \setminus j}, x)$  is in natural exponential family form:

$$p(\theta_j|\theta_{i \in [J] \setminus j}, x) = \exp(\tilde{\eta}_j^T \theta_j - A_j(\tilde{\eta}_j)) \quad (2)$$

with local natural parameter  $\tilde{\eta}_j$  and local log partition function  $A_j$ . Here,  $\tilde{\eta}_j$  may be a function of  $\theta_{i \in [J] \setminus j}$  and  $x$ . If the exponential family assumption above holds for every index  $j$ , then we can write  $\tilde{\eta}_j$  as a sum of products of components of each  $\theta_k$  vector:

$$\tilde{\eta}_j = \sum_{r \in R_j} G_r \prod_{k \in [J] \setminus j} \theta_{kr_k}, \quad (3)$$

where  $G_r$  is a  $D_j$ -sized column vector and  $\theta_{kr_k}$  is a scalar. Here,  $r$  is a vector of length  $J - 1$ . Each entry  $r_k$  of  $r$  is either  $\emptyset$  or an index in  $[D_k]$ . If  $r_k = \emptyset$ , then  $\theta_{kr_k} = 1$ ; otherwise,  $\theta_{kr_k}$  is the  $r_k$ th element of the vector  $\theta_k$ . This notation scheme guarantees that each product contains at most one factor from the vector  $\theta_k$  for each index  $k$ . In particular, the log likelihood is linear in every vector  $\theta_j$ . This property of the log likelihood is guaranteed by Eq. (2). Appendix A.1 contains further details and a proof of Eq. (3).

It follows from Eqs. (1), (2), (3), and the assumed factorization of  $q^*$  that  $\log q_j^*(\theta_j)$  has the form

$$\left( \sum_{r \in R_j} G_r \prod_{k \in [J] \setminus j} \mathbb{E}_{q_r^*}[\theta_{kr_k}] \right)^T \theta_j + C. \quad (4)$$

We see that  $q_j^*$  is in the same exponential family form as  $p(\theta_j | \theta_{i \in [J] \setminus j}, x)$ . Let  $\eta_j$  denote the natural parameter of  $q_j^*$ , and denote the mean parameter of  $q_j^*$  as  $m_j := \mathbb{E}_{q_j^*} \theta_j$ . We see from Eq. (4) that

$$\eta_j = \sum_{r \in R_j} G_r \prod_{k \in [J] \setminus j} m_{kr_k}.$$

Since  $m_j$  is a function of  $\eta_j$ , we have the fixed point equations  $m_j := M_j(m_{i \in [J] \setminus j})$  for mappings  $M_j$  across  $j$  and

$$m := M(m)$$

for the vector of mappings  $M$ .

### 3 Linear response covariance estimation

Let  $V$  denote the covariance matrix of  $\theta$  under the factorized variational distribution  $q^*(\theta)$ , and let  $\Sigma$  denote the covariance matrix of  $\theta$  under the true distribution,  $p(\theta|x)$ :

$$V := \text{Cov}_{q^*} \theta, \quad \Sigma := \text{Cov}_p \theta.$$

$V$  may be a poor estimate of  $\Sigma$ , even when  $m \approx \mathbb{E}_p \theta$ , i.e. when the marginal means match well [4–8]. Our goal is to use the MFVB solution and the techniques of linear response theory [9–11] to construct an improved estimate for  $\Sigma$ .

Define  $p_t(\theta|x)$  such that its log is a linear perturbation of the log posterior:

$$\log p_t(\theta|x) = \log p(\theta|x) + t^T \theta - C(t), \quad (5)$$

where  $C(t)$  is a constant in  $\theta$ . If we assume that  $p(\theta|x)$  is a probability distribution with natural parameters in the interior of the feasible space, then  $p_t(\theta|x)$  is a probability distribution for any  $t$  in an open ball around 0. Since  $C(t)$  normalizes the  $p_t(\theta|x)$  distribution, it is in fact the cumulant generating function of  $p(\theta|x)$ . Further, every (perturbed) conditional distribution  $p_t(\theta_j | \theta_{i \in [J] \setminus j}, x)$  is in the same exponential family as every (unperturbed) conditional distribution  $p(\theta_j | \theta_{i \in [J] \setminus j}, x)$  by construction. So, for each  $t$ , we have mean field variational approximation  $q_t^*$  with marginal means  $m_{t,j} := \mathbb{E}_{q_t^*} \theta_j$  and fixed point equations  $m_{t,j} = M_{t,j}(m_{t,i \in [J] \setminus j})$  across  $j$ . Thus,  $m_t = M_t(m_t)$  as in Section 2. Taking derivatives of the latter relationship with respect to  $t$ , we find

$$\frac{dm_t}{dt^T} = \frac{\partial M_t}{\partial m_t^T} \frac{dm_t}{dt^T} + \frac{\partial M_t}{\partial t^T}. \quad (6)$$

In particular, note that  $t$  is a vector of size  $D$  (the total dimension of  $\theta$ ), and  $\frac{dm_t}{dt^T}$ , e.g., is a matrix of size  $D \times D$  with  $(a, b)$ th entry equal to the scalar  $dm_{t,a}/dt_b$ .

Since  $q_t^*$  is the MFVB approximation for the perturbed posterior  $p_t(\theta|x)$ , we may hope that  $m_t = E_{q_t^*}\theta$  is close to the perturbed-posterior mean  $\mathbb{E}_{p_t}\theta$ . The practical success of MFVB relies on the fact that this approximation is often good in practice. To derive interpretations of the individual terms in Eq. (6), we assume that this equality of means holds, but we indicate where we use this assumption with an approximation sign:  $m_t \approx \mathbb{E}_{p_t}\theta$ . The full derivations of the following equations are given in Appendix A.2.

$$\frac{dm_t}{dt^T} \approx \frac{d}{dt^T} \mathbb{E}_{p_t}\theta = \Sigma_t \quad (7)$$

$$\frac{\partial M_t}{\partial t^T} = \frac{\partial}{\partial t^T} \mathbb{E}_{q_t^*}\theta = V_t \quad (8)$$

$$\frac{dM_t}{dm_t^T} = V_t H_t, \quad (9)$$

where  $\Sigma_t$  is the covariance matrix of  $\theta$  under  $p_t$ ,  $V_t$  is the covariance matrix of  $\theta$  under  $q_t^*$ , and

$$H_t := E_{q^*} \left( \frac{\partial^2 \log p_t(\theta|x)}{\partial \theta \partial \theta^T} \right)$$

Then substituting Eqs. (7), (8), and (9) into Eq. (6), evaluating at  $t = 0$ , and writing  $H$  for  $H_0$  and  $V$  for  $V_0$ , we find

$$\begin{aligned} \hat{\Sigma} &:= \left. \frac{dm_t}{dt^T} \right|_{t=0} \approx \Sigma \\ \hat{\Sigma} &= V H \hat{\Sigma} + V \Rightarrow \\ \hat{\Sigma} &= (I - V H)^{-1} V \end{aligned} \quad (10)$$

Thus, we call  $\hat{\Sigma}$  the LRVB estimate of the true posterior covariance  $\Sigma$ .

### 3.1 Exactness of multivariate normal and SEM

Consider approximating a multivariate normal posterior distribution  $p(\theta|x)$  with MFVB. This case arises, for instance, given a multivariate normal likelihood with fixed covariance  $S$  and an improper uniform prior on the mean parameter  $\mu$ :

$$p(x|\mu) = \prod_{n=1:N} \mathcal{N}(x_n|\mu, S) \text{ and } q^*(\mu) = \prod_{j=1:J} q_j^*(\mu_j)$$

Here,  $\mathcal{N}$  represents the multivariate normal distribution, and the total dimension  $D$  of  $\mu$  is equal to the number of components  $J$ . So  $\mu$  is a  $J$ -length vector for  $J > 1$  with elements  $\mu_1, \dots, \mu_J$ , and  $S$  is a known  $J \times J$  positive definite matrix. Our variational approximation,  $q^*$ , is given by the factorized distribution over mean components.

In this case, it is well known that the MFVB posterior means are correct, but the marginal variances are underestimated if  $S$  is not diagonal. This fact is often used to illustrate the shortcomings of MFVB [4–7].

However, since the posterior means are correctly estimated, the LRVB approximation in Eq. (10) is in fact an equality. That is, for the posterior location of a multivariate normal with known covariance, Eq. (10) is not an approximation, and  $\hat{\Sigma} = \frac{dm_t}{dt^T} = \Sigma$  exactly. A detailed proof of this fact can be found in Appendix B.

This result draws a connection between LRVB and the “supplemented expectation-maximization” (SEM) method of [13]. SEM is an asymptotically exact covariance correction for the EM algorithm that transforms the full-data Fisher information matrix into the observed-data Fisher information matrix using a correction that is formally similar to Eq. (10). In this sense, SEM is a frequentist perspective on a special case of the LRVB correction when the amount of data goes to infinity. More details can be found in Appendix B.

## 4 Scaling

Eq. (10) requires the inverse of a matrix as large as all the unknown natural parameters in the posterior  $p(\theta|x)$ , which generally includes both main parameters and nuisance parameters. In many applications, the number of nuisance parameters may be very large. For example, in the finite mixture of normals model below (Section 6), there is an indicator variable for the cluster assignment for each data point. If we treat these variables as nuisance parameters, the number of nuisance parameters grows with the number of data points  $N$ . As a result, directly computing the matrix inverse in Eq. (10) may be impractical.

However, since the variational covariance  $V$  is block diagonal and  $H$  is often sparse, one may be able to use Schur complements to efficiently calculate sub-matrices of  $\hat{\Sigma}$ . Suppose that our full parameter space,  $\theta$ , can be divided into a small number of variables of primary interest, called  $\alpha$ , and a large (and possibly growing) number of nuisance variables,  $z$ :

$$\theta = \begin{pmatrix} \alpha \\ z \end{pmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_\alpha & \Sigma_{\alpha z} \\ \Sigma_{z\alpha} & \Sigma_z \end{bmatrix}.$$

We can similarly define partitions for  $H$  and  $V$ . Assume also that we have the usual mean field factorization of the variational approximation:  $q^*(\alpha, z) = q^*(\alpha)q^*(z)$ , so that  $V_{\alpha z} = 0$ . (The variational distributions may factor further as well.) We calculate the Schur complement of  $\hat{\Sigma}$  in Eq. (10) with respect to its  $z$ th component to find that

$$\hat{\Sigma}_\alpha = \begin{aligned} & (I_\alpha - V_\alpha H_\alpha - V_\alpha H_{\alpha z} (I_z - V_z H_z)^{-1} V_z H_{z\alpha})^{-1} V_\alpha \end{aligned} \quad (11)$$

Here,  $I_\alpha$  and  $I_z$  refer to  $\alpha$ - and  $z$ -sized identity matrices, respectively. A detailed derivation can be found in Appendix A.2. In cases where  $(I_z - V_z H_z)^{-1}$  can be efficiently calculated, Eq. (11) involves only an  $\alpha$ -sized inverse. A finite mixture of Gaussians model, which we describe in Section 6, is one such case.

## 5 Influence scores

Influence scores are a powerful tool from classical linear regression that describe how much influence a particular data point has on a modeled outcome. They can be used, for example, to identify outliers and investigate the robustness of the linear model [14, 15]. Analogously, it can be useful to know how much Bayesian posterior means depend on the values of individual data points. A number of authors have proposed methods to measure the sensitivity of the posterior distribution to perturbations or deletions of data points both in linear models [16, 17] and more generally [18–20]. LRVB gives a convenient formula to analytically calculate the influence of individual data points as covariances between the  $\theta$  vector and infinitesimal noise added to the data.

Consider the conditional expectation of a single parameter value,  $\theta_i$ , as a function of a single data point,  $x_n$ . Specifically, for notational convenience, define

$$m_{\theta_i}(x_n) = \mathbb{E}_p[\theta_i | x_1, \dots, x_n, \dots, x_N]$$

One measure of the sensitivity of  $\theta_i$  to  $x_n$  is the derivative of this function,  $\frac{d}{dx_n} m_{\theta_i}(x_n) = m'_{\theta_i}(x_n)$ . We will refer to this derivative as an *influence score* for Bayesian models.

To draw a connection between this influence score and covariances, imagine that our observations,  $x$ , are in fact slightly noisy versions of the true data,  $x^*$ . Specifically, our model becomes

$$p(x|x^*, \theta) = p(x|x^*)p(x^*|\theta).$$

In this new model,  $x^*$  are unknown parameters, like  $\theta$ . We assume our posterior beliefs about the true  $x^*$  obey the following assumptions<sup>1</sup>:

$$\begin{aligned} \mathbb{E}(x^*|x) &= x \\ \text{Cov}(x^*|x) &= \Sigma_x \\ S_x &:= \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \Sigma_x \\ 0 \neq S_x &< \infty \\ \text{Higher moments} &= O(\epsilon^p), \text{ for } p > 2 \end{aligned} \tag{12}$$

That is, the covariance matrix  $\Sigma_x$  is proportional to  $\epsilon$ . Conditional on  $x$ ,  $m_{\theta_i}(x_i^*)$  is a random variable that varies as the posterior belief about  $x_i^*$  varies around  $x_i$ . By forming a Taylor expansion of  $m'_{\theta_i}(x_n^*)$  around  $x_n$  we show for any data point  $x_n$  and any parameter  $\theta_i$  that:

$$m'_{\theta_i}(x_n) = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \text{Cov}(\theta_i, x_n^* | x) \tag{13}$$

(Appendix D contains further details.) That is, the limiting value of this covariance as  $\epsilon \rightarrow 0$  can be used to estimate the influence of observation  $x_n$  on the mixture parameters in the spirit of classical linear model influence scores from the statistics literature.

Note that the covariances on the right hand side of Eq. (13) are impossible to compute in naive MFVB, since they involve correlations between distinct mean field components, and difficult to compute using MCMC, since they require estimating a large number of very small covariances with a finite number of draws. However, LRVB leads to a straightforward analytic expression for these covariances.

To derive the LRVB influence scores, we now assume that the parameter space of the posterior can be divided into three types of variables. We have main and nuisance parameters, called  $\alpha$  and  $z$  respectively as in Section 4, and now also  $x^*$ , the unobserved data. As before, we also assume that each has its own variational distribution, i.e.  $q^*(\theta) = q^*(\alpha)q^*(z)q^*(x^*)$ . (The variational distributions may factor still further.) We can write:

$$\theta = \begin{pmatrix} \alpha \\ x^* \\ z \end{pmatrix} \quad \text{and} \quad \Sigma = \begin{bmatrix} \Sigma_\alpha & \Sigma_{\alpha x^*} & \Sigma_{\alpha z} \\ \Sigma_{x^* \alpha} & \Sigma_{x^*} & \Sigma_{x^* z} \\ \Sigma_{z \alpha} & \Sigma_{z x^*} & \Sigma_z \end{bmatrix}.$$

<sup>1</sup> Note that if each observation  $x_n^*$  has only one sufficient statistic, the perturbations can be treated as independent, and  $S_x$  will be the identity. However, if each observation  $x_n^*$  has a vector of sufficient statistics,  $S_x$  must take that structure into account. For example, if  $x_n$  is drawn from a normal distribution centered at  $x_n^*$ , it will have sufficient statistics  $x_n$  and  $x_n^2$ , which will be correlated with one another. These correlations will cause  $S_x$  to be different from the identity in general.

We use a similar partition for  $V$  and  $H$ . Recall that  $\Sigma_x$  is the result of an infinitesimal perturbation and nearly zero, so we express our results in terms of  $S_x$  in Eq. (12). Let  $\Sigma_\alpha$  denote the ordinary LRVB covariance of  $\alpha$  from Eq. (11). The covariance between  $\alpha$  and the infinitesimally perturbed  $x$  then has the following formula:

$$\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \Sigma_{\alpha x^*} = \Sigma_\alpha^{-1} (V_\alpha H_{\alpha x^*} + V_\alpha H_{\alpha z} (I_z - V_z H_z)^{-1} V_z H_{zx^*}) S_x \quad (14)$$

This formula follows from the Schur inverse and taking  $\epsilon \rightarrow 0$ . Details can be found in Appendix D.

## 6 Experiments

Mixture models constitute some of the most popular models for MFVB application [1, 2] and are often used as an example of where MFVB covariance estimates may go awry [5, 6]. We thus illustrate the efficacy of LRVB on the problem of approximating the posterior when the likelihood is a finite mixture of multivariate Gaussians.

### 6.1 Model

We consider a  $K$ -component mixture of  $P$ -dimensional multivariate normals with unknown component means, covariances, and weights. In what follows, the weight  $\pi_k$  is the probability of the  $k$ th component,  $\mathcal{N}$  denotes the multivariate normal distribution,  $\mu_k$  is the  $P$ -dimensional mean of the  $k$ th component, and  $\Lambda_k$  is the  $P \times P$  precision matrix of the  $k$ th component (so  $\Sigma_k := \Lambda_k^{-1}$  is the covariance).  $N$  is the number of data points, and  $x_n$  is the  $n$ th observed  $P$ -dimensional data point. We employ the standard trick of augmenting the data generating process with the latent indicator variables  $z_{nk}$ , where  $n = 1, \dots, N$  and  $k = 1, \dots, K$ , and

$$\begin{aligned} P(z_{nk} = 1) &= \pi_k \\ z_{nk} = 1 &\Rightarrow x_n \sim \mathcal{N}(\mu_k, \Lambda_k^{-1}) \end{aligned}$$

The full likelihood under this augmentation is

$$p(x|\pi, \mu, \Lambda, z) = \prod_{n=1:N} \prod_{k=1:K} \mathcal{N}(x_n|\mu_k, \Lambda_k^{-1})^{z_{nk}} \quad (15)$$

We assign independent variational factors to  $\mu$ ,  $\pi$ ,  $\Lambda$ , and  $z$ .<sup>2</sup> The  $z$  variables are nuisance parameters.

Our goal is to estimate the covariance matrix of the parameters  $\log(\pi)$ ,  $\mu$ ,  $\Lambda$  in the posterior distribution  $p(\pi, \mu, \Lambda|x)$  and to estimate the influence of each data point  $x_n$  on the posterior means of  $\log(\pi)$ ,  $\mu$ ,  $\Lambda$  using LRVB (see Sections 3 and 5).

In addition to the standard MFVB covariance matrices, we will compare the accuracy and speed of our estimates to Gibbs sampling on the augmented model (Eq. (15)) using the function `rmixGibbs` from the R package `bayesm`. Our LRVB implementation relied heavily on linear algebra routines in `RcppEigen` [21]. We evaluate our results both on simulated data and on the MNIST data set [12].

<sup>2</sup> Unlike Section 3.1, the variational posteriors for  $\mu$  factor across components but not within components. That is, for each  $k$ ,  $q^*(\mu_k)$  is a multivariate (not a univariate) normal distribution.

## 6.2 MNIST data set

For a real-world example, we applied LRVB to the unsupervised classification of two digits from the MNIST dataset of handwritten digits. We first preprocess the MNIST dataset by performing principle component analysis on the training data’s centered pixel intensities and keeping the top 25 components. For evaluation, the test data is projected onto the same 25-dimensional subspace found using the training data.

We then treat the problem of separating handwritten 0s from 1s as an unsupervised clustering problem. We limit the dataset to instances labeled as 0 or 1, resulting in 12665 training and 2115 test points. We fit the training data as a mixture of multivariate Gaussians. Here,  $K = 2$ ,  $P = 25$ , and  $N = 12665$ . Then, keeping the  $\mu$ ,  $\Lambda$ , and  $\pi$  parameters fixed, we calculate the expectations of the latent variables  $z$  in Eq. (15) for the test set. We assign test set data point  $x_n$  to whichever component has maximum a posteriori expectation. We count successful classifications as test set points that match their cluster’s majority label and errors as test set points that are different from their cluster’s majority label. By this measure, our test set error rate was 0.08. We stress that we intend only to demonstrate the feasibility of LRVB on a large, real-world dataset rather than to propose practical methods for modeling MNIST.

## 6.3 Covariance experiments

In this section, we check the covariances estimated with Eq. (10) against a Gibbs sampler, which we treat as the ground truth.<sup>3</sup>

For simulations, we generated  $N = 10000$  data points from  $K = 2$  multivariate normal components in  $P = 2$  dimensions. MFVB is expected to underestimate the marginal variance of  $\mu$ ,  $\Lambda$ , and  $\log(\pi)$  when the components overlap since that induces correlation in the posteriors due to the uncertain classification of points between the clusters. These correlations are in violation of the MFVB assumption and cause the MFVB posterior variances to be mis-estimated.

We performed 68 simulations, each of which had at least 500 effective Gibbs samples in each variable—calculated with the R tool `effectiveSize` from the `coda` package [22]. We note that for each of the parameters  $\log(\pi)$ ,  $\mu$ , and  $\Lambda$ , both MH and MFVB produce posterior means close to the ground truth MCMC values, so our key assumption in the LRVB derivations of Section 3 appears to hold.

Each point in Fig. (1) represents the a single parameter in a single simulation. For example, each point on the  $\Lambda$  graph represents the marginal standard deviation of a particular component of the  $\Lambda$  matrix for both the Gibbs sample and an alternative method. The first three graphs show the diagonal standard deviations, and the final graph shows the off-diagonal covariances. Note that the final graph excludes the MFVB estimates since most of the values are zero.

Fig. (1) shows that the raw MFVB covariance estimates are often quite different from the Gibbs sampler results, while the LRVB estimates match the Gibbs sampler closely. Although not shown, the results on the MNIST dataset were as good.

In these simulations, on average LRVB took only 3.40 seconds, whereas the Gibbs sampler took 306.97 seconds. We explore these timing tradeoffs in more detail in Section 6.4.

<sup>3</sup> The likelihood described in Section 6.1 is symmetric under relabeling. When the component locations and shapes have a real-life interpretation, the researcher is generally interested in the uncertainty of  $\mu$ ,  $\Lambda$ , and  $\pi$  for a particular labeling, not the marginal uncertainty over all possible re-labelings. This poses a problem for standard MCMC methods, and we restrict our simulations to regimes where label switching did not occur in our Gibbs sampler. The MFVB solution conveniently avoids this problem since the mean field assumption prevents it from representing more than one mode of the joint posterior.



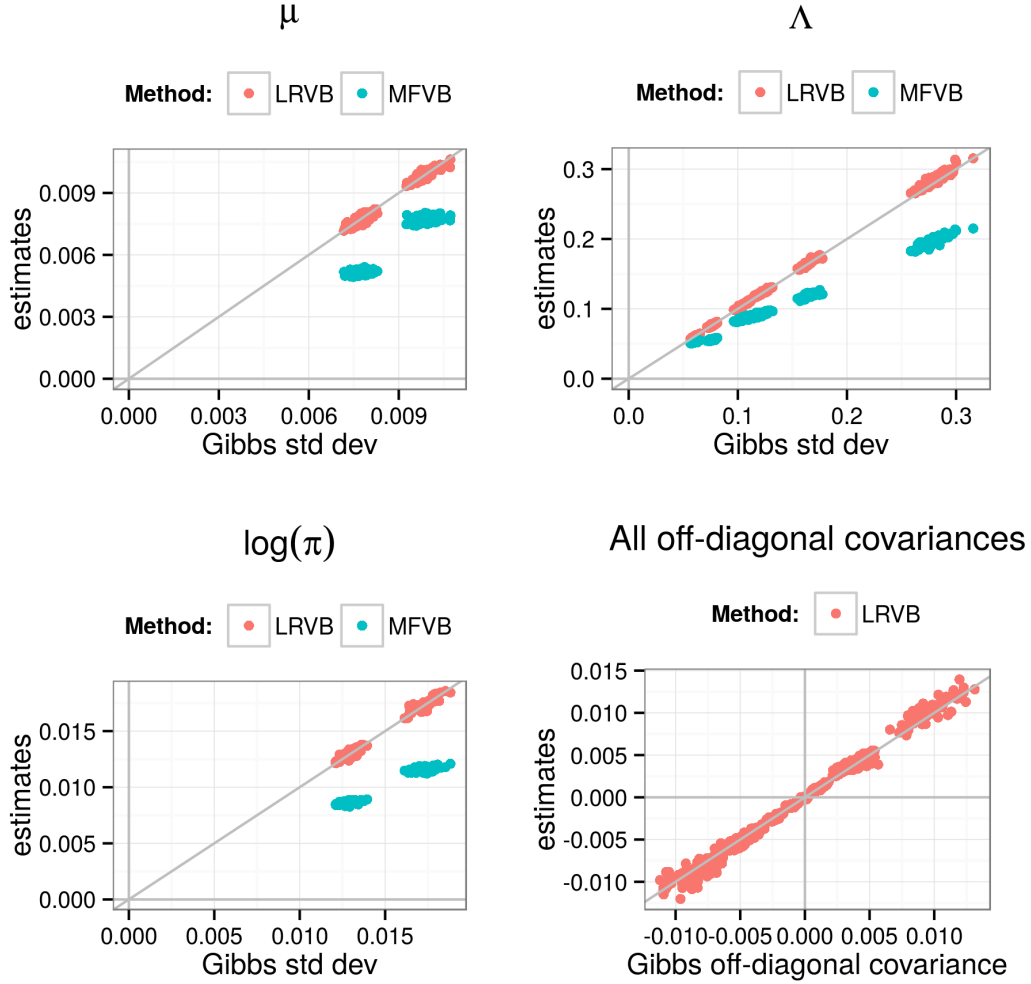


Figure 1: Comparison of estimates of the posterior covariance matrix on simulation data for each model parameter from Gibbs, MFVB, and LRVB methods. In the simulations,  $N = 10000$  (data points),  $K = 2$  (components) and  $P = 2$  (dimensions).

## 6.4 Scaling experiments

In this section we show that, for the finite mixture of multivariate Gaussians model, Eq. (10) scales linearly with  $N$  and polynomially in  $K$  and  $P$ . We also use simulations to experimentally compare the scaling of LRVB running times to Gibbs sampling estimates. We show that LRVB is much faster than Gibbs for the range of parameters we simulated, though Gibbs may be preferable for very high-dimensional problems.

In the terms of Section 4,  $\alpha$  includes the sufficient statistics from  $\mu$ ,  $\pi$ , and  $\Lambda$ , and grows as  $O(KP^2)$ . The sufficient statistics for the variational posterior of  $\mu$  contain the  $P$ -length vectors  $\mu_k$ , for each  $k$ , and the  $(P+1)P/2$  second-order products in the covariance matrix  $\mu_k \mu_k^T$ . Similarly, for

each  $k$ , the variational posterior of  $\Lambda$  involves the  $(P + 1)P/2$  sufficient statistics in the symmetric matrix  $\Lambda_k$  as well as the term  $\log |\Lambda_k|$ . The sufficient statistics for the posterior of  $\pi_k$  are the  $K$  terms  $\log \pi_k$ .<sup>4</sup> This means that, minimally, Eq. (10) will require the inverse of a matrix of size  $O(KP^2)$ .

The sufficient statistics for  $z$  have dimension  $K \times N$ . In other words, the number of nuisance parameters grows with the number of data points, but  $H_z = 0$  for the multivariate normal (Appendix C contains further details), so we can apply Eq. (11) to replace the inverse of an  $O(KN)$ -sized matrix with multiplication by an  $O(KN)$ -sized matrix. Here,  $z$  conveniently corresponds directly to the  $z$  in Section 4.

Since a matrix inverse is cubic in the size of the matrix, the worst-case scaling for LRVB is then  $O(K^2)$  in  $K$ ,  $O(P^6)$  in  $P$  and  $O(N)$  in  $N$ .

In our simulations, shown in Fig. (2), we can see that, in practice, LRVB scales linearly in  $N$  and slightly less than quadratically in  $P$ , which is much better than the theoretical worst case. Note that the vertical axis, the time to run the algorithm, is on the log scale. At every value of  $P$ ,  $K$ , and  $N$  examined here, calculating LRVB is much faster than Gibbs sampling.<sup>5</sup>

## 6.5 Influence score experiments

We now demonstrate the accuracy of the influence score formula, Eq. (14), on simulated data and on the MNIST dataset. We first compare Eq. (14) to numeric derivatives. We then look at some patterns in the data that are made visible by having easy-to-calculate influence scores.

### 6.5.1 Comparison with numeric derivatives

In order to verify that LRVB influence scores are correct for this model, we manually perturbed each component of the data and re-fit to find the new MFVB optimum. In other words, we numerically compute the derivative in Eq. (13).

Our simulation used  $N = 10000$ ,  $P = 2$ , and  $K = 2$ . This was small enough to calculate the influence score for every data point in every dimension. To select sample data points for MNIST influence scores, we selected 5 representative data points with different ranges of  $\mathbb{E}_{q^*} z_n$  values so that some had their posterior probability concentrated in only one component, and some that were uncertainly classified between the two components. We then computed the LRVB influence scores from Eq. (14). For comparison, we again manually perturbed each dimension of each data point and re-optimized.

On MNIST, not counting the time to compute the initial LRVB covariance, calculating the LRVB influence scores took 37 seconds, and the process of perturbing and re-optimizing took 20.7 minutes.

The comparison between numeric differentiation and LRVB influence scores for both a simulation (left) and MNIST (right) is shown in Fig. (3). The influence scores obtained from the two approaches are practically indistinguishable. Though not shown, in both simulations and on MNIST, for each  $\mu$ ,  $\Lambda$ , and  $\log(\pi)$  parameter, the two methods were as similar to one another as in Fig. (3).

<sup>4</sup> Since  $\sum_{k=1:K} \pi_k = 1$ , using  $K$  sufficient statistics involves one redundant parameter. However, this does not violate any of the necessary assumptions for Eq. (10), and it considerably simplifies the calculations. Note that though the perturbation argument of Section 3 requires the natural parameters of  $p(\theta|x)$  to be in the interior of the feasible space, it does not require that the natural parameters of  $p(x|\theta)$  be interior.

<sup>5</sup> For numeric stability we started the optimization procedures for MFVB at the true values, so the time to compute the optimum in our simulations was very fast and not representative of practice. On real data, the optimization time will depend on the quality of the starting point. Consequently, the times shown for LRVB are only the times to compute the LRVB estimate. The optimization times were on the same order. The Gibbs sampling time was linearly rescaled to the amount of time necessary to achieve 1000 effective samples in the slowest-mixing component of any parameter.

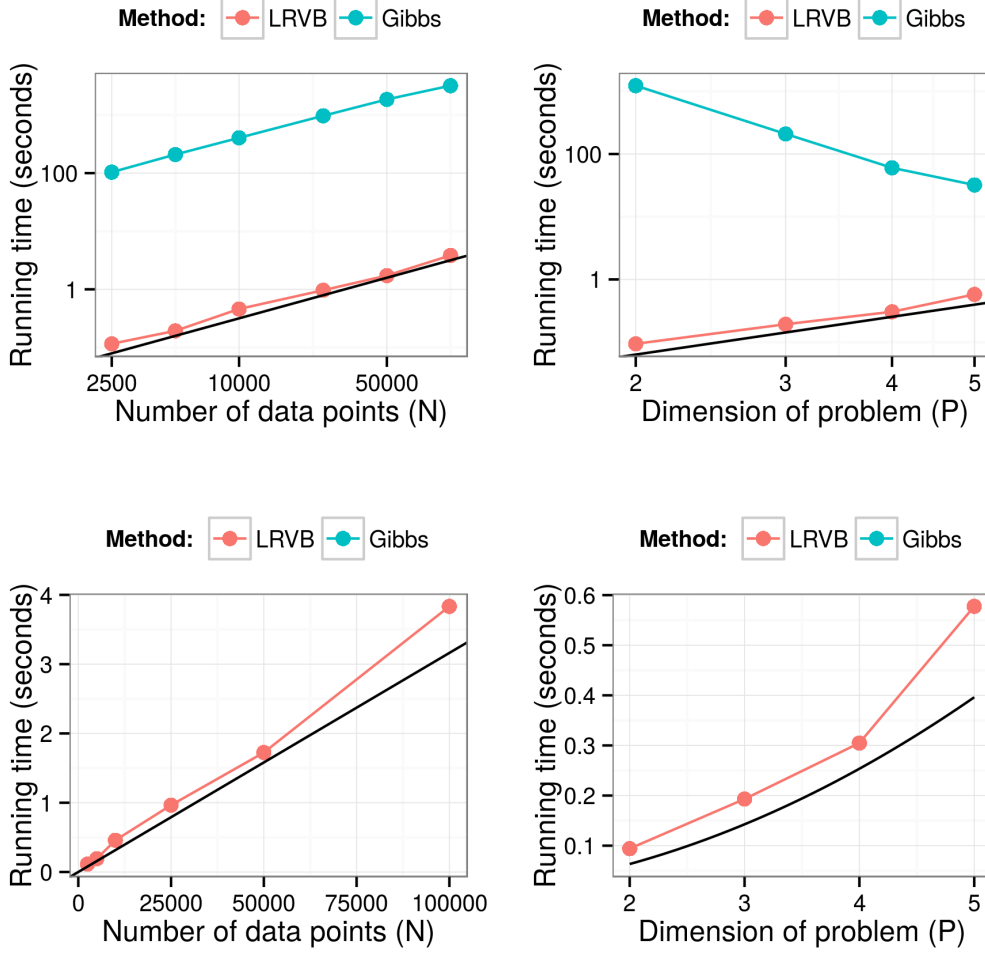


Figure 2: Scaling of LRVB and Gibbs on simulation data in both log and linear scales. Before taking logs, the line in the (N) graph is  $y \propto x$ , and in the (P) graph, it is  $y \propto x^2$ .

### 6.5.2 Influence score data

One can see interesting patterns in the data with influence scores. Consider, for example, the simulated data, which is depicted in Fig. (4) and has moderately overlapping components. The graph shows the effect on  $\mu_{11}$  of perturbing the  $x_{n1}$  (horizontal) coordinate of each datapoint. One can see that the component's mean is essentially determined by the points that are assigned to it. Interestingly, points on the border between the two components reverse the sign of their effect. This is caused by changes in  $\mathbb{E}_{q^*} z_n$  that more than counterbalance their effects on  $\mathbb{E}_{q^*} \mu$ .

As can be seen in the simulated data of Fig. (4), the situation becomes more complex when the components overlap even more. Data points that are distant from a component center have nearly as much influence as data points well within the component.

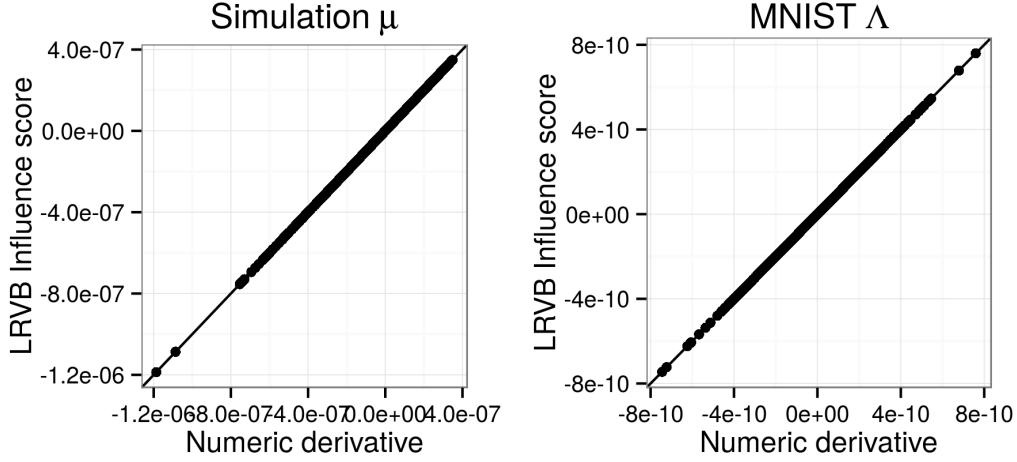


Figure 3: LRVB influence score accuracy. The simulated data uses all  $x_n$ , and the MNIST data uses five representative  $x_n$ .

Finally, we consider influence scores in the MNIST data set. We selected 100 data points with a 0 or 1 label uniformly at random. In Fig. (4), we considered the influence of one particular dimension of each data point on one particular dimension of each component mean. Recall from Section 6.2 that each data point and component mean is 25-dimensional. Now we wish to derive a single influence score for the effect of each data point vector-valued  $x_n$  on each vector-valued component mean  $\mu_k$ .

To that end, we define the influence of a data point  $x_n$  on  $\mu_k$  as the directional derivative of  $\|\mu_k\|_2^2$  with respect to  $x_n$ . That is, we calculate the vector  $\partial\|\mu_k\|_2^2/\partial x_n$  using LRVB as described above. Then we compute

$$\text{Influence of } x_n \text{ on } \mu_k := \max_{\delta: \|\delta\|=1} \left( \frac{\partial\|\mu_k\|_2^2}{\partial x_n^T} \delta \right)$$

The resulting influences are plotted in Fig. (5). Each point corresponds to a data point  $x_n$ . Each sub-figure corresponds to a different component mean parameter. The horizontal axis value for point  $x_n$  is the logit of  $\mathbb{E}_{q^*} z_{nk}$  (capped at  $\pm 15$ ), which measures the posterior probability that  $x_n$  came from component  $k$ .

The two components show very different patterns. Component 0, the mode with mostly hand-written zeroes, has much higher influence amongst points that are classified within it than component 1.

## 7 Conclusion

The lack of accurate covariance estimates from the widely used mean-field variational Bayes (MFVB) methodology has been a longstanding shortcoming of MFVB. We have demonstrated that our method, linear response variational Bayes (LRVB), augments MFVB to deliver these covariance estimates in time that scales linearly with the number of data points. We have also shown how to use LRVB

to quickly calculate influence scores, a measure of the influence of each data point on posterior parameter means. Our experiments have focused on mixtures of multivariate Gaussians since these have traditionally been used to illustrate the difficulties with MFVB covariance estimation. We hope that in future work our results can be extended to more complex models, including latent Dirichlet allocation and Bayesian nonparametric models, where MFVB has proven its practical success.

## References

- [1] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.
- [2] D. M. Blei and M. I. Jordan. Variational inference for Dirichlet process mixtures. *Bayesian Analysis*, 1(1):121–143, 2006.
- [3] M. D. Hoffman, D. M. Blei, C. Wang, and J. Paisley. Stochastic variational inference. *Journal of Machine Learning Research*, 14(1):1303–1347, 2013.
- [4] D. J. C. MacKay. *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press, 2003. Chapter 33.
- [5] C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer, New York, 2006. Chapter 10.
- [6] R. E. Turner and M. Sahani. Two problems with variational expectation maximisation for time-series models. In D. Barber, A. T. Cemgil, and S. Chiappa, editors, *Bayesian Time Series Models*. 2011.
- [7] B. Wang and M. Titterton. Inadequacy of interval estimates corresponding to variational Bayesian approximations. In *Workshop on Artificial Intelligence and Statistics*, pages 373–380, 2004.
- [8] H. Rue, S. Martino, and N. Chopin. Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations. *Journal of the Royal Statistical Society: Series B (statistical methodology)*, 71(2):319–392, 2009.
- [9] H. J. Kappen and F. B. Rodriguez. Efficient learning in Boltzmann machines using linear response theory. *Neural Computation*, 10(5):1137–1156, 1998.
- [10] M. Opper and O. Winther. Variational linear response. In *Advances in Neural Information Processing Systems*, 2003.
- [11] M. Welling and Y. W. Teh. Linear response algorithms for approximate inference in graphical models. *Neural Computation*, 16(1):197–221, 2004.
- [12] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [13] X. L. Meng and D. B. Rubin. Using EM to obtain asymptotic variance-covariance matrices: The SEM algorithm. *Journal of the American Statistical Association*, 86(416):899–909, 1991.
- [14] S. Chatterjee and A. S. Hadi. Influential observations, high leverage points, and outliers in linear regression. *Statistical Science*, pages 379–393, 1986.

- [15] R. D. Cook. Assessment of local influence. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 133–169, 1986.
- [16] I. Guttman and D. Peña. A Bayesian look at diagnostics in the univariate linear model, 1992. Universidad Carlos III de Madrid: Working Paper 92-21.
- [17] D. Peña and I. Guttman. Comparing probabilistic methods for outlier detection in linear models. *Biometrika*, 80(3):603–610, 1993.
- [18] B. P. Carlin and N. G. Polson. An expected utility approach to influence diagnostics. *Journal of the American Statistical Association*, 86(416):1013–1021, 1991.
- [19] F. Peng and D. K. Dey. Bayesian analysis of outlier problems using divergence measures. *Canadian Journal of Statistics*, 23(2):199–213, 1995.
- [20] H. Zhu, J. G. Ibrahim, and N. Tang. Bayesian influence analysis: a geometric approach. *Biometrika*, 98(2):307–323, 2011.
- [21] D. Bates and D. Eddelbuettel. Fast and elegant numerical linear algebra using the RcppEigen package. *Journal of Statistical Software*, 52(5):1–24, 2013.
- [22] Martyn Plummer, Nicky Best, Kate Cowles, and Karen Vines. Coda: Convergence diagnosis and output analysis for mcmc. *R News*, 6(1):7–11, 2006.
- [23] R. W. Keener. *Theoretical Statistics*. Springer, 2010.

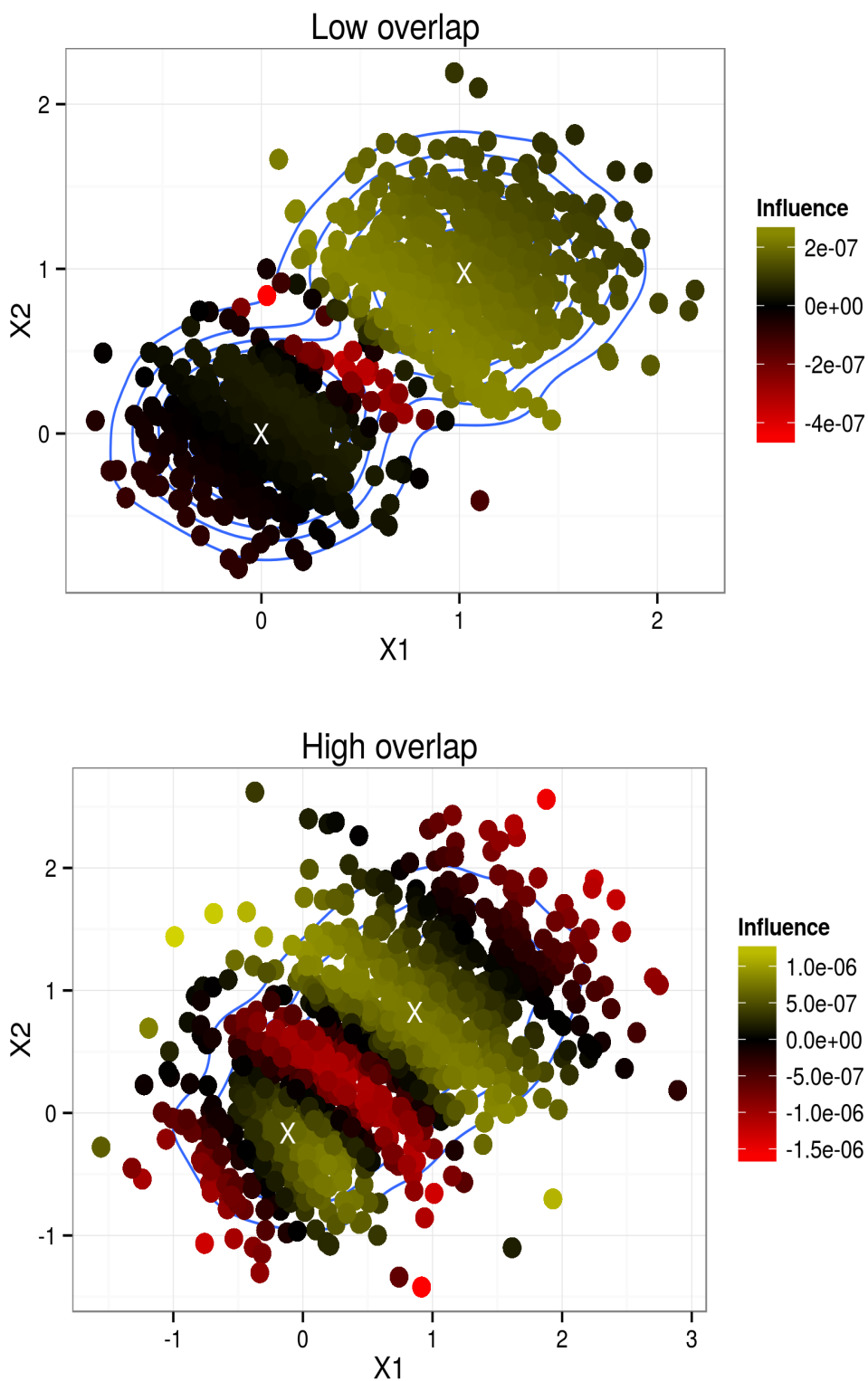


Figure 4: Influence scores for components with different amounts of overlap. Each graph shows the influence of  $x_{n1}$  on  $\mu_{11}$ , which is mean of the upper-right hand component. (X) indicates a component posterior mean.

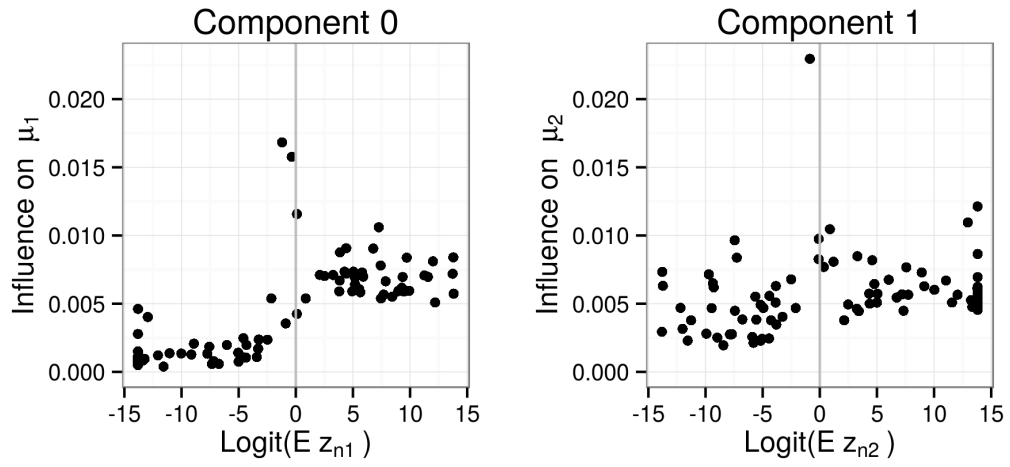


Figure 5: Different influence patterns for the two Gaussians in the MNIST dataset. The 100 data points were chosen randomly. The vertical axis shows the maximum directional derivative of  $\mu_k$  with respect to changes in the data point, and the horizontal axis shows the (capped) logit posterior probability that the point came from that component.



# Supplementary Material

## A Derivations

### A.1 MFVB for conditional exponential families

First, we require some notation for indexing  $\theta$ . Recall that the MFVB assumption partitions the components of  $\theta$  into  $J$  groups according to the factorization

$$q(\theta) = \prod_{j=1}^J q_j(\theta_j) = \prod_{j=1}^J q(\theta_j)$$

We follow the common abuse of notation of defining the variational functions through their arguments by writing  $q(\theta_j)$  for  $q_j(\theta_j)$ .

Each  $\theta_j$  is be a  $D_j$ -dimensional vector, with  $\sum_{j=1}^J D_j = D$ , where the whole  $\theta$  vector has dimension  $D$ . Here and below, we will define the set  $[J] := \{1, \dots, J\}$ .

Let  $R(\theta)$  denote a vector of length  $\prod_{j \in [J]} (D_j + 1)$  that is the vector of all possible products of the form  $\prod_{j \in [J]} \theta_{ji_j}$ , where for each  $j$ ,  $i_j \in \{\emptyset, 1, \dots, D_j\}$ . That is,  $R(\theta)$  is the vector of all possible products of terms of  $\theta$  where with at most one term from each  $\theta_j$ , and we define  $\theta_{j\emptyset} := 1$ . Let  $R_j(\theta)$  denote the same vector, but excluding terms from  $\theta_j$ , and  $R_{jj'}(\theta)$  denote the same vector, but excluding both  $\theta_j$  and  $\theta_{j'}$ .

To aid in intuition, it will sometimes be useful to explicitly write the inner product of  $R(\theta)$  with a vector  $G$  as a sum of products of components of  $\theta_j$ . Let  $G$  be a  $|R(\theta)|$ -length vector with element  $G_i$ , and let  $R(\theta)_i$  denote the  $i$ th row of  $R(\theta)$ . Then define

$$G^T R(\theta) = \sum_{i=1}^{|R(\theta)|} G_i R(\theta)_i := \sum_{r \in R} G_r \prod_{j \in [J]} \theta_{jr_j} \quad (16)$$

Here, we have “overloaded” the definition of  $R$  to express the sum over different products of  $\theta$ . We define  $r_j$  as the index of  $\theta_j$  in the row of  $R(\theta)$  corresponding to  $r_j$ ,  $G_r$  as the element of  $G$  corresponding to that index set, and the sum over  $r \in R$  as the sum over all rows.

We define the inner product of  $G$  with  $R_j(\theta)$  similarly, only with the result as a  $D_j$ -length vector. Specifically, define

$$\sum_{r \in R_j} G_r \prod_{k \in [J] \setminus j} \theta_{kr_k}$$

such that  $G_r$  is a  $D_j$ -sized column vector. Finally, define

$$\sum_{r \in R_{jj'}} G_r \prod_{k \in [J] \setminus \{j, j'\}} \theta_{kr_k}$$

the same way, except where  $G_r$  is a  $D_j \times D_{j'}$  matrix.

This notation is intended to make it easy to express sums of products of elements of  $\theta$  in such a way that no two terms from a single  $\theta_j$  are multiplied together. The value of this notation will hopefully become clear in the following lemmas.

**Lemma A.1.** Suppose Eq. (2) holds across all  $j$ ; that is,

$$p(\theta_j | \theta_{i \in [J] \setminus j}, x) = \exp(\tilde{\eta}_j^T \theta_j - A_j(\tilde{\eta}_j)).$$

Then the posterior  $p(\theta|x)$  can be written in the form

$$\log p(\theta|x) = \sum_{r \in R} G_r \prod_{j \in [J]} \theta_{jr_j} + C \quad (17)$$

where the terms  $G_r$  and  $C$  are constant in all  $\theta$ <sup>6</sup>.

*Proof.* We see that  $\log p(\theta|x) = \log p(\theta_j | \theta_{i \in [J] \setminus j}, x) + \log p(\theta_{i \in [J] \setminus j} | x)$  depends on  $\theta_j$  only via the first term in the sum. By Eq. (2),

$$p(\theta_j | \theta_{i \in [J] \setminus j}, x) = \exp(\tilde{\eta}_j^T \theta_j - A_j(\tilde{\eta}_j))$$

It follows that  $\log p(\theta|x)$  is linear in the vector  $\theta_j$ . But this is true for all  $\theta_j$ , and the above form for  $p(\theta|x)$  follows.  $\square$

**Lemma A.2.** Suppose Eq. (2) holds across all  $j$ . Then, for the natural parameter  $\tilde{\eta}_j$ , we have the following equations:

$$\begin{aligned} \tilde{\eta}_j &= \sum_{r \in R_j} G_r \prod_{k \in [J] \setminus j} \theta_{kr_k} \\ \tilde{\eta}_j &= \frac{\partial \log p(\theta|x)}{\partial \theta_j} \\ H &= \frac{\partial \eta}{\partial m^T} = E_{q^*} \left( \frac{\partial^2 \log p(\theta|x)}{\partial \theta \partial \theta^T} \right) \end{aligned} \quad (18)$$

Here, each  $G_r$  is an  $D_j$ -length vector that is constant with respect to  $\theta$ .

*Proof.* The first result follows by collecting the terms for the  $i$ th component of  $\theta_j$  in Eq. (17) and applying Bayes' theorem. The second is simply observing that differentiating with respect to  $\theta_j$  is a notationally tidy way to collect the  $j$  terms.

For the third result, recall from Eq. (4) that

$$\begin{aligned} \eta_j &= \mathbb{E}_{q^*}[\tilde{\eta}_j] \\ &= \mathbb{E}_{q^*} \left[ \sum_{r \in R_j} G_r \prod_{k \in [J] \setminus j} \theta_{kr_k} \right] \\ &= \sum_{r \in R_j} G_r \prod_{k \in [J] \setminus j} m_{kr_k} \Rightarrow \\ \frac{\partial \eta_j}{\partial m_{j'}^T} &= \sum_{r \in R_{j'}} G'_r \prod_{k \in [J] \setminus \{j, j'\}} m_{kr_k} \end{aligned}$$

---

<sup>6</sup>Strictly speaking,  $C$  is redundant since  $\{\emptyset, \dots, \emptyset\} \in R$ . Here and below, for additional clarity we will always write a constant.

$$\begin{aligned}
&= \mathbb{E}_{q^*} \left[ \sum_{r \in R_{jj'}} G'_r \prod_{l \in [J] \setminus \{j, j'\}} \theta_{k_r k_l} \right] \\
&= \mathbb{E}_{q^*} \left[ \frac{\partial \log p(\theta|x)}{\partial \theta_i \partial \theta_{j'}^T} \right]
\end{aligned}$$

Note that this proof relied on the fact that only one element of  $\theta_{j'}$  is in each product term, which allowed us to exchange the derivative with respect to the expectation with the expectation of the derivative with respect to  $\theta_{j'}$ .  $\square$

## A.2 Linear response

We here derive the three equalities in Eqs. (7), (8), and (9), which appear respectively as three propositions below. In these propositions, we assume that  $p(\theta|x)$  is in the exponential family as above. We will further assume that all natural parameters (for  $p$  or variational approximations) are in the interior of the parameter space. The fact that the natural parameters are on the interior of the feasible space means that there exists an open ball around them that is also feasible. Let that ball have radius  $\delta$ , and let  $t$  be within a  $\delta$  ball of the origin. Then  $p_t(\theta|x)$  is well defined for  $t$  in an open set containing zero. These assumptions will allow us to apply dominated convergence (cf. Section 2.3 of [23]).

**Proposition A.3.**  $\frac{d}{dt} \mathbb{E}_{p_t} \theta = \Sigma_t$ .

*Proof.*

$$\begin{aligned}
\frac{d}{dt} \mathbb{E}_{p_t} \theta &= \frac{d}{dt} \int_{\theta} \theta e^{t^T \theta - c(t)} p(\theta|x) d\theta \quad \text{by the definition of } p_t \text{ in Eq. (5)} \\
&= \int_{\theta} \theta \left[ \frac{d}{dt} e^{t^T \theta - c(t)} \right] p(\theta|x) d\theta \quad \text{by dominated convergence} \\
&= \int_{\theta} \theta \theta^T e^{t^T \theta - c(t)} p(\theta|x) d\theta - \int_{\theta} \theta e^{t^T \theta - c(t)} p(\theta|x) d\theta \cdot \frac{dc(t)}{dt} \\
&= \mathbb{E}_{p_t} [\theta \theta^T] - \mathbb{E}_{p_t} [\theta] \mathbb{E}_{p_t} [\theta]^T = \Sigma_t
\end{aligned}$$

$\square$

To approximate  $\frac{dm_t}{dt}$ , we assume not only that  $m_t \approx \mathbb{E}_{p_t} \theta$  for any particular  $t$  but further that  $m_t$  tracks the true mean  $\mathbb{E}_{p_t} \theta$  as  $t$  varies. In this case, by Proposition A.3, we have

$$\frac{dm_t}{dt} \approx \frac{d}{dt} \mathbb{E}_{p_t} \theta = \Sigma_t,$$

the first (approximate) equality in Eq. (7).

To derive the final two equalities in Eqs. (8) and (9), we make use of the following lemma.

**Lemma A.4.**  $M_{t,j}$  depends on  $t$  only via  $\eta_{t,j}$ , the natural parameter of the  $q_{t,j}^*$  distribution. And  $\frac{dM_{t,j}}{d\eta_{t,j}^T} = \Sigma_{q_{t,j}^*}$ .

*Proof.* The first part of the lemma follows from writing the definition of  $M_{t,j}$ :

$$M_{t,j} = \mathbb{E}_{q_{t,j}^*} \theta_j = \int_{\theta_j} \theta_j \exp(\eta_{t,j}^T \theta_j - A_j(\eta_{t,j})) d\theta_j.$$

For the second part,

$$\begin{aligned} \frac{dM_{t,j}}{d\eta_{t,j}^T} &= \int_{\theta_j} \frac{d}{d\eta_{t,j}^T} \theta_j \exp(\eta_{t,j}^T \theta_j - A_j(\eta_{t,j})) d\theta_j \quad \text{by dominated convergence} \\ &= \int_{\theta_j} \theta_j [\theta_j^T - \mathbb{E}_{q_{t,j}^*} \theta_j^T] \exp(\eta_{t,j}^T \theta_j - A_j(\eta_{t,j})) d\theta_j \\ &= \Sigma_{q_{t,j}^*} \end{aligned}$$

□

**Proposition A.5.**  $\frac{\partial M_t}{\partial t^T} = V_t$ .

*Proof.* By Lemma A.4, we have for any indices  $i$  and  $j$  in  $[J]$  that

$$\frac{\partial M_{t,j}}{\partial t_i^T} = \frac{dM_{t,j}}{d\eta_{t,j}^T} \frac{\partial \eta_{t,j}}{\partial t_i^T}, \quad (19)$$

where the first factor is also given by Lemma A.4. It remains to find the second factor,  $\frac{\partial \eta_{t,j}}{\partial t_i^T}$ . By the discussion after Eq. (2) and the construction of  $p_t$ , the natural parameter  $\tilde{\eta}_{t,j}$  of  $p_t(\theta_j | \theta_{i \in [J] \setminus j}, x)$  satisfies

$$\tilde{\eta}_{t,j} = \sum_{r \in R_j} G_r \prod_{k \in [J] \setminus j} \theta_{kr_k} + t_j.$$

So, as in the derivation of Eq. (4), the natural parameter  $\eta_{t,j}$  of  $q_j^*(\theta_j)$  satisfies

$$\eta_{t,j} i = \sum_{r \in R_j} G_r \prod_{k \in [J] \setminus j} m_{t,kr_k} + t_j \quad (20)$$

for  $m_{t,r} := \mathbb{E}_{q_{t,r}^*} \theta_r$ .

Let  $d_j$  be the dimension of  $\theta_j$  and hence the dimension of  $\eta_{t,j}$  and  $t_j$ . Hence,

$$\frac{\partial \eta_{t,j}}{\partial t_i^T} = \begin{cases} I_{d_j} & j = i \\ 0_{d_j, d_i} & \text{else} \end{cases},$$

where  $I_a$  is the identity matrix of dimension  $a$ , and  $0_{a,b}$  is the all zeros matrix of dimension  $a \times b$ .

Finally, by Eq. (19), Lemma A.4, and the expression for  $\frac{\partial \eta_{t,j}}{\partial t_i^T}$  just obtained, we have

$$\frac{\partial M_t}{\partial t^T} = V_t I_D = V_t.$$

□

**Proposition A.6.**  $\frac{dM_t}{dm_{t,i}^T} = V_t \frac{\partial \eta_t}{\partial m_{t,i}^T}$ .

*Proof.* By Lemma A.4 and analogous to Eq. (19), we have

$$\frac{\partial M_{t,j}}{\partial m_{t,i}^T} = \frac{dM_{t,j}}{d\eta_{t,j}^T} \frac{\partial \eta_{t,j}}{\partial m_{t,i}^T}. \quad (21)$$

The result follows immediately from Lemma A.4.

□

## B Multivariate normal posteriors and SEM

For any target distribution  $p(\theta|x)$ , it is well-known that MFVB cannot be used to estimate the covariances between the components of  $\theta$ . In particular, if  $q^*$  is the estimate of  $p(\theta|x)$  returned by MFVB,  $q^*$  will have a block-diagonal covariance matrix—no matter the form of the covariance of  $p(\theta|x)$ . By contrast, the next result shows that the LRVB covariance estimate is exactly correct in the case where the target distribution,  $p(\theta|x)$ , is (multivariate) normal.

In order to prove this result, we will rely on the following lemma.

**Lemma B.1.** *Consider a target posterior distribution characterized by  $p(\theta|x) = \mathcal{N}(\theta|\mu, \Sigma)$ , where  $\mu$  and  $\Sigma$  may depend on  $x$ , and  $\Sigma$  is invertible. Let  $\theta = (\theta_1, \dots, \theta_J)$ , and consider a MFVB approximation to  $p(\theta|x)$  that factorizes as  $q(\theta) = \prod_j q(\theta_j)$ . Then the variational posterior means are the true posterior means; i.e.  $m_j = \mu_j$  for all  $j$  between 1 and  $J$ .*

*Proof.* The derivation of MFVB for the multivariate normal can be found in Section 10.1.2 of [5]; we highlight some key results here. Let  $\Lambda = \Sigma^{-1}$ . Let the  $j$  index on a row or column correspond to  $\theta_j$ , and let the  $-j$  index correspond to  $\{\theta_i : i \in [J] \setminus j\}$ . E.g., for  $j = 1$ ,

$$\Lambda = \begin{bmatrix} \Lambda_{11} & \Lambda_{1,-1} \\ \Lambda_{-1,1} & \Lambda_{-1,-1} \end{bmatrix}.$$

By the assumption that  $p(\theta|x) = \mathcal{N}(\theta|\mu, \Sigma)$ , we have

$$\log p(\theta_j | \theta_{i \in [J] \setminus j}, x) = -\frac{1}{2}(\theta_j - \mu_j)^T \Lambda_{jj}(\theta_j - \mu_j) + (\theta_j - \mu_j)^T \Lambda_{j,-j}(\theta_{-j} - \mu_{-j}) + C, \quad (22)$$

where the final term is constant with respect to  $\theta_j$ . It follows that

$$\begin{aligned} \log q_j^*(\theta_j) &= \mathbb{E}_{q_i^* : i \in [J] \setminus j} \log p(\theta, x) + C \\ &= -\frac{1}{2}\theta_j^T \Lambda_{jj} \theta_j + \theta_j \mu_j \Lambda_{jj} - \theta_j \Lambda_{j,-j} (\mathbb{E}_{q^*} \theta_{-j} - \mu_{-j}). \end{aligned}$$

So

$$q_j^*(\theta_j) = \mathcal{N}(\theta_j | m_j, \Lambda_{jj}^{-1}),$$

with mean parameters

$$m_j = \mathbb{E}_{q_j^*} \theta_j = \mu_j - \Lambda_{jj}^{-1} \Lambda_{j,-j} (m_{-j} - \mu_{-j}) \quad (23)$$

as well as an equation for  $\mathbb{E}_{q^*} \theta^T \theta$ .

Note that  $\Lambda_{jj}$  must be invertible, for if it were not,  $\Sigma$  would not be invertible.

The solution  $m = \mu$  is a unique stable point for Eq. (23), since the fixed point equations for each  $j$  can be stacked and rearranged to give

$$\begin{aligned} m - \mu &= - \begin{bmatrix} 0 & \Lambda_{11}^{-1} \Lambda_{12} & \cdots & \Lambda_{11}^{-1} \Lambda_{1(J-1)} & \Lambda_{11}^{-1} \Lambda_{1J} \\ \vdots & & \ddots & & \vdots \\ \Lambda_{JJ}^{-1} \Lambda_{J1} & \Lambda_{JJ}^{-1} \Lambda_{J2} & \cdots & \Lambda_{JJ}^{-1} \Lambda_{J(J-1)} & 0 \end{bmatrix} (m - \mu) \\ &= - \begin{bmatrix} \Lambda_{11}^{-1} & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & & & \vdots \\ 0 & & \ddots & & 0 \\ \vdots & & & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & \Lambda_{JJ}^{-1} \end{bmatrix} \begin{bmatrix} 0 & \Lambda_{12} & \cdots & \Lambda_{1(J-1)} & \Lambda_{1J} \\ \vdots & & \ddots & & \vdots \\ \Lambda_{J1} & \Lambda_{J2} & \cdots & \Lambda_{J(J-1)} & 0 \end{bmatrix} (m - \mu) \Leftrightarrow \end{aligned}$$

$$\begin{aligned}
0 &= \begin{bmatrix} \Lambda_{11} & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & & & \vdots \\ 0 & & \ddots & & 0 \\ \vdots & & & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & \Lambda_{JJ} \end{bmatrix} (m - \mu) + \\
&\quad \begin{bmatrix} 0 & \Lambda_{12} & \cdots & \Lambda_{1(J-1)} & \Lambda_{1J} \\ \vdots & & \ddots & & \vdots \\ \Lambda_{J1} & \Lambda_{J2} & \cdots & \Lambda_{J(J-1)} & 0 \end{bmatrix} (m - \mu) \Leftrightarrow \\
0 &= \Lambda (m - \mu) \Leftrightarrow \\
m &= \mu.
\end{aligned}$$

The last step follows from the assumption that  $\Sigma$  (and hence  $\Lambda$ ) is invertible. It follows that  $\mu$  is the unique stable point of Eq. (23).  $\square$

**Proposition B.2.** *Assume we are in the setting of Lemma B.1, where additionally  $\mu$  and  $\Sigma$  are on the interior of the feasible parameter space. Then the LRVB covariance estimate exactly captures the true covariance,  $\hat{\Sigma} = \Sigma$ .*

*Proof.* Consider the perturbation for LRVB defined in Eq. (5). By perturbing the log likelihood, we change both the true means  $\mu_t$  and the variational solutions,  $m_t$ . The result is a valid density function since the original  $\mu$  and  $\Sigma$  are on the interior of the parameter space. By Lemma B.1, the MFVB solutions are exactly the true means, so  $m_{t,j} = \mu_{t,j}$ , and the derivatives are the same as well. This means that the first term in Eq. (10) is not approximate, i.e.

$$\frac{dm_t}{dt^T} = \frac{d}{dt^T} \mathbb{E}_{p_t} \theta = \Sigma_t,$$

It follows from the arguments in Appendix B that the LRVB covariance matrix is exact, and  $\hat{\Sigma} = \Sigma$ .  $\square$

## B.1 Comparison with supplemented expectation-maximization

This result about the multivariate normal distribution draws a connection between LRVB corrections and the “supplemented expectation-maximization” (SEM) method of [13]. SEM is an asymptotically exact covariance correction for the EM algorithm that transforms the full-data Fisher information matrix into the observed-data Fisher information matrix using a correction that is formally similar to Eq. (10). In this section, we argue that this similarity is not a coincidence; in fact the SEM correction is an asymptotic version of LRVB with two variational blocks, one for the missing data and one for the unknown parameters.

Although LRVB as described here requires a prior (unlike SEM, which supplements the MLE), the two covariance corrections coincide when the full information likelihood is approximately log quadratic and proportional to the posterior,  $p(\theta|x)$ . This might be expected to occur when we have a large number of independent data points informing each parameter—i.e., when a central limit theorem applies and the priors do not affect the posterior. In the full information likelihood, some terms may be viewed as missing data, whereas in the Bayesian model the same terms may be viewed as latent parameters, but this does not prevent us from formally comparing the two methods.

We can draw a term-by-term analogy with the equations in [13]. We denote variables from the SEM paper with a superscript “SEM” to avoid confusion. MFVB does not differentiate between missing data and parameters to be estimated, so our  $\theta$  corresponds to  $(\theta^{SEM}, Y_{mis}^{SEM})$  in [13]. SEM is an asymptotic theory, so we may assume that  $(\theta^{SEM}, Y_{mis}^{SEM})$  have a multivariate normal distribution, and that we are interested in the mean and covariance of  $\theta^{SEM}$ .

In the E-step of [13], we replace  $Y_{mis}^{SEM}$  with its conditional expectation given the data and other  $\theta^{SEM}$ . This corresponds precisely to Eq. (23), taking  $\theta_j = Y_{mis}^{SEM}$ . In the M-step, we find the maximum of the log likelihood with respect to  $\theta^{SEM}$ , keeping  $Y_{mis}^{SEM}$  fixed at its expectation. Since the mode of a multivariate normal distribution is also its mean, this, too, corresponds to Eq. (23), now taking  $\theta_j = \theta^{SEM}$ .

It follows that the MFVB and EM fixed point equations are the same; i.e., our  $M$  is the same as their  $M^{SEM}$ , and our  $\partial M / \partial m$  of Eq. (21) corresponds to the transpose of their  $DM^{SEM}$ , defined in Eq. (2.2.1) of [13]. Since the “complete information” corresponds to the variance of  $\theta^{SEM}$  with fixed values for  $Y_{OBS}^{SEM}$ , this is the same as our  $\Sigma_{q^*, 11}$ , the variational covariance, whose inverse is  $I_{oc}^{-1}$ . Taken all together, this means that equation (2.4.6) of [13] can be re-written as our Eq. (10).

$$\begin{aligned} V^{SEM} &= I_{oc}^{-1} (I - DM^{SEM})^{-1} \Rightarrow \\ \Sigma &= V \left( I - \left( \frac{\partial M}{\partial m^T} \right)^T \right)^{-1} = \left( I - \frac{\partial M}{\partial m^T} \right)^{-1} V \end{aligned}$$

## C Multivariate normal mixture details

In this section we derive the basic formulas needed to calculate Eq. (10) and Eq. (14) for a finite mixture of normals, which is the model used in Section 6. We will follow the notation introduced in Section 6.1.

Let each observation,  $x_n$ , be a  $P \times 1$  vector. We will denote the  $P$ th component of the  $n$ th observation  $x_n$ , with a similar pattern for  $z$  and  $\mu$ . We will denote the  $p, q$ th entry in the matrix  $\Lambda_k$  as  $\Lambda_{k,pq}$ . The data generating process is as follows:

$$\begin{aligned} \log P(x_n | z_n, \mu, \Lambda) &= \sum_{n=1}^N z_{nk} \log \phi_k(x_n) + C \\ \log \phi_k(x) &= -\frac{1}{2} (x - \mu_k)^T \Lambda_k (x - \mu_k) + \frac{1}{2} \log |\Lambda_k| + C \\ \log P(z_{nk} | \pi_k) &= \sum_{k=1}^K z_{nk} \log \pi_k + C \\ \log P(z, \mu, \pi, \Lambda | x) &= \sum_{n=1}^N \sum_{k=1}^K z_{nk} \left( \log \pi_k - \frac{1}{2} (x_n - \mu_k)^T \Lambda_k (x_n - \mu_k) + \frac{1}{2} \log |\Lambda_k| \right) + C \end{aligned}$$

In all our results we simply used improper, flat priors, though it would be trivial to incorporate conjugate priors.

From the assumptions in Eq. (12), the posterior expectation of  $x^*$  will always have  $\mathbb{E}_q x^* = x$ , so for notational convenience we can simply drop the  $*$  and apply the LRVB formulas as if  $x$  were a random parameter. However, it is useful to remember that the parameter  $x$  is different from the variable we condition on – we are actually estimating  $p(\alpha, z, x^* | x)$ .

The parameters  $\mu_k$ ,  $\Lambda_k$ ,  $\pi$ , and  $z_n$  will each be given their own variational distribution. By standard results, the variational distributions will be:

$$\begin{aligned} q_{\mu_k} &= \text{Multivariate Normal} \\ q_{\Lambda_k} &= \text{Wishart} \\ q_{\pi} &= \text{Dirichlet} \\ q_{z_n} &= \text{Multinoulli (one multinomial draw)} \\ q_{x_n^*} &= \text{Multivariate Normal} \end{aligned}$$

The sufficient statistics for  $\mu_k$  are all terms of the form  $\mu_{kp}$  and  $\mu_{kp}\mu_{kq}$ . Consequently, the sub-vector of  $\theta$  corresponding to  $\mu_k$  is

$$\theta_{\mu_k} = \begin{pmatrix} \mu_{k1} \\ \vdots \\ \mu_{kp} \\ \mu_{k1}\mu_{k1} \\ \mu_{k1}\mu_{k2} \\ \vdots \\ \mu_{kP}\mu_{kP} \end{pmatrix}$$

We will only save one copy of  $\mu_{kp}\mu_{kq}$  and  $\mu_{kq}\mu_{kp}$ , so  $\theta_{\mu_k}$  has length  $P + \frac{1}{2}(P+1)P$ . For all the parameters, we denote the complete stacked vector without a  $k$  subscript:

$$\theta_{\mu} = \begin{pmatrix} \theta_{\mu_1} \\ \vdots \\ \theta_{\mu_K} \end{pmatrix}$$

The sufficient statistics for  $x^*$  are analogous to those for  $\mu$ .

The sufficient statistics for  $\Lambda_k$  are all the terms  $\Lambda_{k,pq}$  and the term  $\log |\Lambda_k|$ . Again, since  $\Lambda$  is symmetric, we do not keep redundant terms, so  $\theta_{\Lambda_k}$  has length  $1 + \frac{1}{2}(P+1)P$ .

The sufficient statistics for  $\pi$  is the  $K$ -vector  $(\log \pi_1, \dots, \log \pi_K)$ .

The sufficient statistics for  $z$  is simply the  $N \times K$  values  $z_{nk}$  themselves.

In terms of Section 5, we have

$$\begin{aligned} \alpha &= \begin{pmatrix} \theta_{\mu} \\ \theta_{\Lambda} \\ \theta_{\pi} \end{pmatrix} \\ z &= (\theta_z) \\ x &= (\theta_x) \end{aligned}$$

That is, we are primarily interested in the covariance of the sufficient statistics of  $\mu$ ,  $\Lambda$ , and  $\pi$ ,  $z$  are nuisance parameters, and  $x^*$  is the “unobserved” data.

To put the log likelihood in terms useful for LRVB, we must express it in terms of the sufficient statistics, taking into account the fact the  $\theta$  vector does not store redundant terms (e.g. it will only keep  $\Lambda_{ab}$  for  $a < b$  since  $\Lambda$  is symmetric).

$$-\frac{1}{2}(x_n - \mu_k)^T \Lambda_k (x_n - \mu_k) = -\frac{1}{2}\text{trace}\left(\Lambda_k (x_n - \mu_k)(x_n - \mu_k)^T\right)$$



$$\begin{aligned}
&= -\frac{1}{2} \sum_a \sum_b (\Lambda_{k,ab} (x_{n,a} - \mu_{k,a}) (x_{n,b} - \mu_{k,b})) \\
&= -\frac{1}{2} \sum_a \sum_b (\Lambda_{k,ab} \mu_{k,a} \mu_{k,b} - \Lambda_{k,ab} x_{n,a} \mu_{k,b} - \Lambda_{k,ab} x_{n,b} \mu_{k,a} + \Lambda_{k,ab} x_{n,a} x_{n,b}) \\
&= -\frac{1}{2} \sum_a \Lambda_{k,aa} (\mu_k^2)^a + \sum_a \Lambda_{k,aa} x_{n,a} \mu_{k,a} - \frac{1}{2} \sum_a \Lambda_{k,aa} (x_n^2)^2 - \\
&\quad \frac{1}{2} \sum_{a \neq b} \Lambda_{k,ab} \mu_{k,a} \mu_{k,b} + \sum_{a \neq b} \Lambda_{k,ab} x_{n,a} \mu_{k,b} - \frac{1}{2} \sum_{a \neq b} \Lambda_{k,ab} x_{n,a} x_{n,b} \\
&= -\frac{1}{2} \sum_a \Lambda_{k,aa} (\mu_k^2)^a + \sum_a \Lambda_{k,aa} x_{n,a} \mu_{k,a} - \frac{1}{2} \sum_a \Lambda_{k,aa} (x_n^2)^2 - \\
&\quad \sum_{a < b} \Lambda_{k,ab} \mu_{k,a} \mu_{k,b} + \sum_{a < b} \Lambda_{k,ab} (x_{n,a} \mu_{k,b} + x_{n,b} \mu_{k,a}) - \sum_{a < b} \Lambda_{k,ab} x_{n,a} x_{n,b}
\end{aligned}$$

The MFVB updates and covariances in  $V$  are all given by properties of standard distributions. To compute the LRVB corrections, it only remains to calculate the hessian of  $H$ . These terms can be read directly off the posterior. First we calculate derivatives with respect to components of  $\mu$ .

$$\begin{aligned}
\frac{\partial^2 H}{\partial \mu_{k,a} \partial \Lambda_{k,ab}} &= \sum_i z_{nk} x_{n,b} \\
\frac{\partial^2 H}{\partial (\mu_{k,a} \mu_{k,b}) \partial \Lambda_{k,ab}} &= -\left(\frac{1}{2}\right)^{1(a=b)} \sum_n z_{nk} \\
\frac{\partial^2 H}{\partial \mu_{k,a} \partial z_{nk}} &= \sum_b \Lambda_{k,ab} x_{n,b} \\
\frac{\partial^2 H}{\partial (\mu_{k,a} \mu_{k,b}) \partial z_{nk}} &= -\left(\frac{1}{2}\right)^{1(a=b)} \Lambda_{k,ab}
\end{aligned}$$

All other  $\mu$  derivatives are zero. For  $\Lambda$ ,

$$\begin{aligned}
\frac{\partial^2 H}{\partial \Lambda_{k,ab} \partial z_{nk}} &= -\left(\frac{1}{2}\right)^{1(a=b)} (x_{n,a} x_{n,b} - \mu_{k,a} x_{n,b} - \mu_{k,b} x_{n,a} + \mu_{k,a} \mu_{k,b}) \\
\frac{\partial^2 H}{\partial \log |\Lambda_k| \partial z_{nk}} &= \frac{1}{2}
\end{aligned}$$

The remaining  $\Lambda$  derivatives are zero. The only nonzero second derivatives for  $\log \pi$  are to  $Z$  and are given by

$$\frac{\partial^2 H}{\partial \log \pi_j \partial z_{nk}} = 1$$

To calculate the influence scores, we additionally need second derivatives involving  $x$ .

$$\frac{\partial^2 H}{\partial x_{n,a} \partial \mu_{k,b}} = z_{nk} \Lambda_{k,ab}$$

$$\begin{aligned}
\frac{\partial^2 H}{\partial x_{n,a} \partial \Lambda_{k,ab}} &= z_{nk} \mu_{k,b} \\
\frac{\partial^2 H}{\partial x_{n,a} \partial z_{nk}} &= \sum_b \mu_{k,b} \Lambda_{k,ab} \\
\frac{\partial^2 H}{\partial (x_{n,a} x_{n,b}) \partial \Lambda_{k,ab}} &= -z_{nk} \left(\frac{1}{2}\right)^{1(a=b)} \\
\frac{\partial^2 H}{\partial (x_{n,a} x_{n,b}) \partial z_{nk}} &= -\Lambda_{k,ab} \left(\frac{1}{2}\right)^{1(a=b)}
\end{aligned}$$

All other second derivatives involving  $x$  are zero. Note in particular that  $H_{zz} = 0$ , allowing efficient calculation of Eq. (11).

## D Influence scores

### D.1 Influence score derivations

In this section, we derive the formulas in Section 5. We will follow the notation there defined. Consider the “influence score” given by the derivative of the conditional expectation:

$$\begin{aligned}
m_{\theta_i}(x_n) &= \mathbb{E}_p[\theta_i | x_1, \dots, x_n, \dots, x_N] \\
\frac{d}{dx_n} m_{\theta_i}(x_n) &:= m'_{\theta_i}(x_n)
\end{aligned}$$

A Taylor expansion of  $m_{\theta_i}(x_n^*)$  around  $x_n$  gives

$$\begin{aligned}
m_{\theta_i}(x_n^*) &= m_{\theta_i}(x_n) + m'_{\theta_i}(x_n)(x_n^* - x_n) + \\
&\quad O\left((x_n^* - x_n)^2\right)
\end{aligned}$$

Multiplying both sides by  $(x_n^* - x_n)$  and taking expectations conditional on  $x$  gives

$$\begin{aligned}
&\mathbb{E}[(m_{\theta_i}(x_n^*) - m_{\theta_i}(x_n))(x_n^* - x_n) | x] \\
&= m'_{\theta_i}(\theta_i, x_n) \mathbb{E}[(x_n^* - x_n)^2 | x] + O\left((x_n^* - x_n)^3\right)
\end{aligned}$$

On the left side,

$$\begin{aligned}
&\mathbb{E}[(m_{\theta_i}(x_n^*) - m_{\theta_i}(x_n))(x_n^* - x_n) | x] \\
&= \mathbb{E}[\mathbb{E}[(m_{\theta_i}(x_n^*) - m_{\theta_i}(x_n))(x_n^* - x_n) | x_n^*] | x] \\
&= \mathbb{E}[\mathbb{E}[(\mathbb{E}[\theta_i | x_1, \dots, x_n^*, \dots, x_N] - m_{\theta_i}(x_n))(x_n^* - x_n) | x_n^*] | x] \\
&= \mathbb{E}[\mathbb{E}[(\theta_i - m_{\theta_i}(x_n))(x_n^* - x_n) | x_n^*] | x] \\
&= \mathbb{E}[(\theta_i - m_{\theta_i}(x_n))(x_n^* - x_n) | x] \\
&= \text{Cov}(\theta_i, x_n^* | x)
\end{aligned}$$

On the right side,

$$m'_{\theta_i}(\theta_i, x_n) \mathbb{E}[(x_n^* - x_n)^2 | x] + O\left((x_n^* - x_n)^3\right)$$

$$= m'_{\theta_i}(\theta_i, x_n) \epsilon + O\left(\epsilon^{\frac{3}{2}}\right)$$

So that

$$m'_{\theta_i}(\theta_i, x_n) = \frac{1}{\epsilon} \text{Cov}(\theta_i, x_n^* | x) + O\left(\epsilon^{\frac{3}{2}}\right)$$

...which is Eq. (13).

We now assume that our parameter space can be divided into types of variables:  $\alpha$  and  $z$  as before, and  $x$ , the perturbed data. As before, we also assume that each has its own variational distribution.

$$\theta = \begin{pmatrix} \alpha \\ x \\ z \end{pmatrix}$$

$$\Sigma = \begin{bmatrix} \Sigma_\alpha & \Sigma_{\alpha x^*} & \Sigma_{\alpha z} \\ \Sigma_{x^* \alpha} & \Sigma_{x^*} & \Sigma_{x^* z} \\ \Sigma_{z \alpha} & \Sigma_{z x^*} & \Sigma_z \end{bmatrix}$$

As before, we use a similar partition for  $V$  and  $H$ . Specifically,

$$V = \begin{bmatrix} V_\alpha & 0 & 0 \\ 0 & V_{x^*} & 0 \\ 0 & 0 & V_z \end{bmatrix}$$

$$H = \begin{bmatrix} H_\alpha & H_{\alpha x^*} & H_{\alpha z} \\ H_{x^* \alpha} & H_{x^*} & H_{x^* z} \\ H_{z \alpha} & H_{z x^*} & H_z \end{bmatrix}$$

We are interested in  $\Sigma_{\alpha x^*}$ , the covariance between  $\alpha$  and  $x$ , which can be interpreted as influence scores.

The matrix  $\Sigma_{x^*}$  is the result of an infinitesimal perturbation, and so will be nearly zero. We will write

$$\Sigma_{x^*} = \epsilon S_{x^*}$$

Note that  $S_{x^*}$  is not necessarily diagonal if the variational distribution for each datapoint  $x_n$  is multidimensional. Applying formula Eq. (11) to eliminate  $z$ , we have

$$\begin{aligned} \begin{bmatrix} \Sigma_\alpha & \Sigma_{\alpha x^*} \\ \Sigma_{x^* \alpha} & \Sigma_{x^*} \end{bmatrix} &= \left[ \begin{pmatrix} I_\alpha - V_\alpha H_\alpha & -V_\alpha H_{\alpha x^*} \\ -\epsilon S_{x^*} H_{x^* \alpha} & I_{x^*} \end{pmatrix} - \begin{pmatrix} V_\alpha H_{\alpha z} \\ \epsilon S_{x^*} H_{x^* z} \end{pmatrix} Q_z \begin{pmatrix} V_z H_{z \alpha} & V_z H_{z x^*} \end{pmatrix} \right]^{-1} \begin{pmatrix} V_\alpha & 0 \\ 0 & \epsilon S_{x^*} \end{pmatrix} \\ &= \left[ \begin{pmatrix} I_\alpha - V_\alpha H_\alpha & -V_\alpha H_{\alpha x^*} \\ -\epsilon S_{x^*} H_{x^* \alpha} & I_{x^*} \end{pmatrix} - \begin{pmatrix} V_\alpha H_{\alpha z} Q_z V_z H_{z \alpha} & V_\alpha H_{\alpha z} Q_z V_z H_{z x^*} \\ \epsilon S_{x^*} H_{x^* z} Q_z V_z H_{z \alpha} & \epsilon S_{x^*} H_{x^* z} Q_z V_z H_{z x^*} \end{pmatrix} \right]^{-1} \begin{pmatrix} V_\alpha & 0 \\ 0 & \epsilon S_{x^*} \end{pmatrix} \\ &= \begin{bmatrix} I_\alpha - V_\alpha H_\alpha - V_\alpha H_{\alpha z} Q_z V_z H_{z \alpha} & -(V_\alpha H_{\alpha x^*} + V_\alpha H_{\alpha z} Q_z V_z H_{z x^*}) \\ -\epsilon (S_{x^*} H_{x^* \alpha} + S_{x^*} H_{x^* z} Q_z V_z H_{z \alpha}) & I_{x^*} - \epsilon S_{x^*} H_{x^* z} Q_z V_z H_{z x^*} \end{bmatrix}^{-1} \begin{pmatrix} V_\alpha & 0 \\ 0 & \epsilon S_{x^*} \end{pmatrix} \\ &= \begin{bmatrix} I_\alpha - V_\alpha H_\alpha - V_\alpha H_{\alpha z} Q_z V_z H_{z \alpha} & -Q_{\alpha x^*} \\ -\epsilon Q_{x^* \alpha} & I_{x^*} - \epsilon Q_{x^*} \end{bmatrix}^{-1} \begin{pmatrix} V_\alpha & 0 \\ 0 & \epsilon S_{x^*} \end{pmatrix} \end{aligned}$$

In the last step we have defined a some placeholder matrices called  $Q$  to simplify subsequent expressions:

$$\begin{aligned} Q_z &:= (I_z - V_z H_z)^{-1} \\ Q_{\alpha x^*} &:= V_\alpha H_{\alpha x^*} + V_\alpha H_{\alpha z} Q_z V_z H_{zx^*} \\ Q_{x^* \alpha} &:= S_{x^*} H_{xq} + S_{x^*} H_{x^* z} Q_z V_z H_{z\alpha} \\ Q_{x^*} &:= S_{x^*} H_{x^* z} Q_z V_z H_{z\alpha} \end{aligned}$$

This may appear to be a complicated expression, but it can be considerably simplified by using the fact that  $\Sigma_{\alpha x^*} \propto \epsilon$  and  $\epsilon \approx 0$ , which allows us to eliminate all  $\epsilon$  terms that are second-order or higher. We can also use the Taylor expansion of the matrix inverse that gives, for  $\epsilon$  small, and invertible matrix  $A$ ,

$$(I - \epsilon B)^{-1} = I + \epsilon B + O(\epsilon^2)$$

Again applying a Schur complement, we can write the expression for the upper-left corner:

$$\begin{aligned} \Sigma_\alpha &= \left( I_\alpha - V_\alpha H_\alpha - V_\alpha H_{\alpha z} Q_z V_z H_{z\alpha} - \epsilon Q_{\alpha x^*} (I_{x^*} - \epsilon Q_{x^*})^{-1} Q_{x^* \alpha} \right)^{-1} V_\alpha \\ &= (I_\alpha - V_\alpha H_\alpha - V_\alpha H_{\alpha z} Q_z V_z H_{z\alpha})^{-1} V_\alpha + O(\epsilon) \end{aligned}$$

Note that as  $\epsilon \rightarrow 0$ , this gives the ordinary LRVB estimate for  $\Sigma_\alpha$ , as expected. Infinitesimal perturbations to our data do not change our beliefs about the posterior covariance. Next, the Schur complement formula for the upper right corner gives

$$\begin{aligned} \Sigma_{\alpha x^*} &= \epsilon \Sigma_\alpha^{-1} \left( Q_{\alpha x^*} (I_{x^*} - \epsilon Q_{x^*})^{-1} \right) S_{x^*} \\ &= \epsilon \Sigma_\alpha^{-1} \left( Q_{\alpha x^*} (I_{x^*} + \epsilon Q_{x^*} + O(\epsilon^2)) \right) S_{x^*} \\ &= \epsilon \Sigma_\alpha^{-1} Q_{\alpha x^*} S_{x^*} + O(\epsilon^2) \\ &= \epsilon \Sigma_\alpha^{-1} \left( V_\alpha H_{\alpha x^*} + V_\alpha H_{\alpha z} (I_z - V_z H_z)^{-1} V_z H_{zx^*} \right) S_{x^*} + O(\epsilon^2) \end{aligned}$$

Taking limits gives Eq. (14).